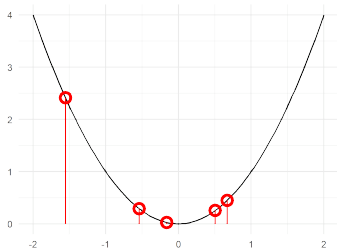


## Regression Losses: L2-loss



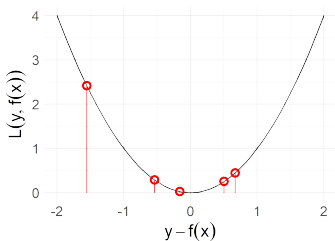
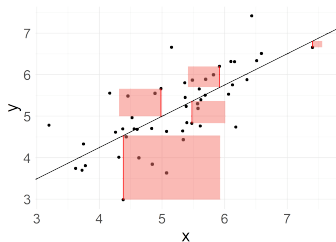
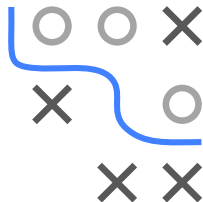
- Derive the risk minimizer of the L2-loss
- Derive the optimal constant model for the L2-loss



# L2-LOSS

$$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2 \quad \text{or} \quad L(y, f(\mathbf{x})) = 0.5 (y - f(\mathbf{x}))^2$$

- Tries to reduce large residuals (if residual is twice as large, loss is 4 times as large), hence outliers in  $y$  can become problematic
- Analytic properties: convex, differentiable (gradient no problem in loss minimization)
- Residuals = Pseudo-residuals:  $\tilde{r} = -\frac{\partial 0.5(y-f(\mathbf{x}))^2}{\partial f(\mathbf{x})} = y - f(\mathbf{x}) = r$



## L2-LOSS: RISK MINIMIZER

Let us consider the (true) risk for  $\mathcal{Y} = \mathbb{R}$  and the L2-Loss

$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$  with unrestricted  $\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathbb{R}^g\}$ .

- By the law of total expectation

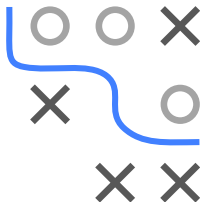
$$\begin{aligned}\mathcal{R}_L(f) &= \mathbb{E}_{xy} [L(y, f(\mathbf{x}))] = \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{y|\mathbf{x}} [L(y, f(\mathbf{x})) \mid \mathbf{x} = \mathbf{x}]] \\ &= \mathbb{E}_{\mathbf{x}} [\mathbb{E}_{y|\mathbf{x}} [(y - f(\mathbf{x}))^2 \mid \mathbf{x} = \mathbf{x}]] .\end{aligned}$$

- Since  $\mathcal{H}$  is unrestricted, at any point  $\mathbf{x} = \mathbf{x}$ , we can predict any value  $c$  we want. The best point-wise prediction is the cond. mean

$$f^*(\mathbf{x}) = \operatorname{argmin}_c \mathbb{E}_{y|\mathbf{x}} [(y - c)^2 \mid \mathbf{x} = \mathbf{x}] \stackrel{(*)}{=} \mathbb{E}_{y|\mathbf{x}} [y \mid \mathbf{x}] .$$

(\*) follows from:

$$\begin{aligned}\operatorname{argmin}_c \mathbb{E} [(y - c)^2] &= \operatorname{argmin}_c \underbrace{\mathbb{E} [(y - c)^2] - (\mathbb{E}[y] - c)^2}_{= \operatorname{Var}[y - c] = \operatorname{Var}[y]} + (\mathbb{E}[y] - c)^2 \\ &= \operatorname{argmin}_c \operatorname{Var}[y] + (\mathbb{E}[y] - c)^2 = \mathbb{E}[y] .\end{aligned}$$

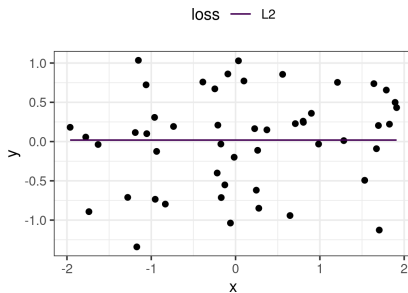


# L2-LOSS: OPTIMAL CONSTANT MODEL

The optimal constant model in terms of the (theoretical) risk for the L2 loss is the expected value over  $y$ :

$$f(\mathbf{x}) = \mathbb{E}_{y \mid \mathbf{x}}[y \mid \mathbf{x}] \stackrel{\text{drop } \mathbf{x}}{=} \mathbb{E}_y[y]$$

The optimizer of the empirical risk is  $\bar{y}$  (the empirical mean over  $y^{(i)}$ ), which is the empirical estimate for  $\mathbb{E}_y[y]$ .



# L2-LOSS: OPTIMAL CONSTANT MODEL / 2

## Proof:

For the optimal constant model  $f(\mathbf{x}) = \theta$  for the L2-loss  $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$  we solve the optimization problem

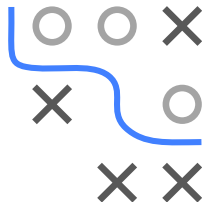
$$\arg \min_{f \in \mathcal{H}} \mathcal{R}_{\text{emp}}(f) = \arg \min_{\theta \in \mathbb{R}} \sum_{i=1}^n (y^{(i)} - \theta)^2.$$

We calculate the first derivative of  $\mathcal{R}_{\text{emp}}$  w.r.t.  $\theta$  and set it to 0:

$$\frac{\partial \mathcal{R}_{\text{emp}}(\theta)}{\partial \theta} = -2 \sum_{i=1}^n (y^{(i)} - \theta) \stackrel{!}{=} 0$$

$$\sum_{i=1}^n y^{(i)} - n\theta = 0$$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n y^{(i)} =: \bar{y}.$$

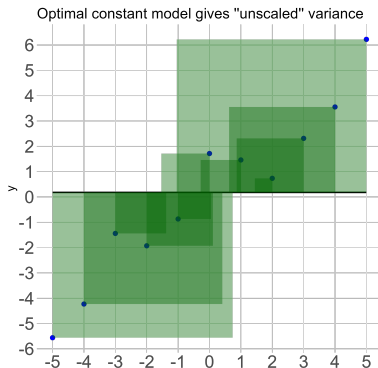


# L2 LOSS MEANS MINIMIZING VARIANCE

Rethinking what we just did: We optimized for a constant whose squared distance to all data points is minimal (in sum, or on average). This turned out to be the mean.

What if we calculate the loss of  $\hat{\theta} = \bar{y}$ ? That's  $\mathcal{R}_{\text{emp}} = \sum_{i=1}^n (y^{(i)} - \bar{y})^2$ .

Average this by  $\frac{1}{n}$  or  $\frac{1}{n-1}$  to obtain variance.



- Generally, if model yields unbiased predictions,  $\mathbb{E}_{y | \mathbf{x}} [y - f(\mathbf{x}) | \mathbf{x}] = 0$ , using  $L_2$ -loss means minimizing variance of model residuals
- Same holds for the pointwise construction / conditional distribution considered before